

## Component Analysis for Computer Vision

Fernando De la Torre



Tutorial European Conference on Computer Vision  
May 2006

## Component Analysis

- Computer Vision & Image Processing
  - Structure from motion.
  - Spectral graph methods for segmentation.
  - Appearance and shape models.
  - Fundamental matrix estimation and calibration.
  - Compression.
  - Classification.
  - Dimensionality reduction and visualization.
- Signal Processing
  - Spectral estimation, system identification (e.g. Kalman filter), sensor array processing (e.g. cocktail problem, eco cancellation), blind source separation, ...
- Computer Graphics
  - Compression (BRDF), synthesis, ...
- Speech, bioinformatics, combinatorial problems.



Component Analysis for Computer Vision

F. De la Torre

ECCV-06

2

## Outline

- Introduction
- Generative models
  - Principal Component Analysis (PCA).
  - Non-negative Matrix Factorization (NMF).
  - Independent Component Analysis (ICA).
- Discriminative models
  - Linear Discriminant Analysis (LDA).
  - Oriented Component Analysis (OCA).
  - Canonical Correlation Analysis (CCA).
  - Relevant Component Analysis (RCA).
- Standard extensions of linear models
  - Latent variable models.
  - Tensor factorization.
  - Kernel methods.



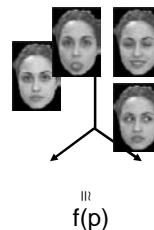
Component Analysis for Computer Vision

F. De la Torre

ECCV-06

3

## Why Subspace Methods?



- Learning: High dimensional data lie in a low dimensional manifold.
- Estimation: Many cv problems (SFM, calibration, ...) can be posed as subspace estimation.
- Better generalization (noise removal).
- Simple parameterization.
- Lower computational complexity.
- Closed form solution and global minimum.



Component Analysis for Computer Vision

F. De la Torre

ECCV-06

4

## Generative Models

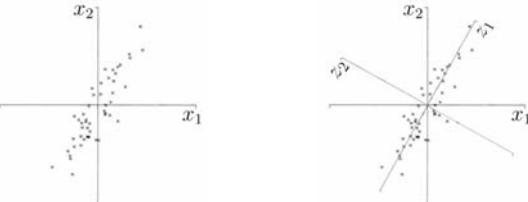
$$\mathbf{D} \approx \mathbf{BC}$$

- Principal Component Analysis/Singular Value Decomposition.
  - 1) Robust PCA/SVD.
  - 2) PCA with uncertainty and missing data.
  - 3) Parameterized PCA.
  - 4) PCA over continuous spaces.
  - 5) Filtered PCA.
  - 6) PCA of rotated images.
  - 7) Mixture of subspaces.
- Non-Negative Matrix Factorization.
  - 8) PCA and NMF for spectral clustering.
- Independent Component Analysis.



## Principal Component Analysis (PCA)

(Pearson, 1901; Hotelling, 1933; Mardia et al., 1979; Jolliffe, 1986; Diamantaras, 1996)



- PCA finds the directions of maximum variation of the data based on linear correlation.
- PCA decorrelates the original variables.



## PCA



$\overset{\text{signal}}{=}$

$$D = \underbrace{\left[ \begin{array}{c} d_1 \\ d_2 \\ \vdots \\ d_n \end{array} \right]}_{n=\text{images}} \approx BC + \mu \mathbf{1}_n^T$$

$D \in \mathbb{R}^{d \times n} \quad B \in \mathbb{R}^{d \times k} \quad C \in \mathbb{R}^{k \times n} \quad \mu \in \mathbb{R}^{d \times 1}$

$$\underset{\text{mean}}{D} \approx \mu \underset{\text{basis}}{\left[ \begin{array}{c} \square \\ \vdots \\ \square \end{array} \right]} + c_1 \underset{\text{component}}{\left[ \begin{array}{c} \square \\ \vdots \\ \square \end{array} \right]} + c_2 \underset{\text{component}}{\left[ \begin{array}{c} \square \\ \vdots \\ \square \end{array} \right]} + \dots + c_k \underset{\text{component}}{\left[ \begin{array}{c} \square \\ \vdots \\ \square \end{array} \right]}$$

Assuming 0 mean data, the basis  $B$  that preserve the maximum variation of the signal is given by eigenvectors  $DD^T$ .

$$d \left| \begin{array}{c} d \\ DD^T B = B \Lambda \end{array} \right.$$



## Snap-shot method & SVD

- If  $d \gg n$  (e.g. images 100\*100 vs. 300 samples) no  $DD^T$ .
- $DD^T$  and  $D^T D$  have the same eigenvalues (energy) and related eigenvectors (by  $D$ ).
- $B$  is a linear combination of the data! (Sirovich, 1987)  
 $DD^T B = B \Lambda \quad B = D \alpha \quad \cancel{D^T D^T D \alpha = D^T D \alpha \Lambda}$
- $[\alpha, \Lambda] = \text{eig}(D^T D) \quad B = D \alpha (\text{diag}(\text{diag}(\Lambda)))^{-0.5}$
- SVD factorizes the data matrix  $D$  as:  
 $DD^T = U \Lambda U^T$   
 $D = U \Sigma V^T \quad D^T D = V \Lambda V^T$   
 $B = U^T \Sigma V^T \quad \Lambda = C C^T$   
 $B^T B = I \quad C C^T = \Lambda$   
 $U^T U = I \quad V^T V = I \quad \Sigma \text{ diagonal}$



## Error function for PCA

- PCA minimizes the reconstruction error.  
(Eckardt & Young, 1936; Gabriel & Zamir, 1979; Baldi & Hornik, 1989; Shum et al., 1995; de la Torre & Black, 2003a)

$$E_1(\mathbf{B}, \mathbf{C}) = \sum_{i=1}^n \|\mathbf{d}_i - \mathbf{B}\mathbf{c}_i\|_2^2 = \|\mathbf{D} - \mathbf{BC}\|_F$$

- Not unique solution:  $\mathbf{BR}^{-1}\mathbf{C} = \mathbf{BC}$   $\mathbf{R} \in \Re^{k \times k}$
- To obtain same PCA solution R has to satisfy:

$$\begin{aligned}\hat{\mathbf{B}} &= \mathbf{BR} & \hat{\mathbf{C}} &= \mathbf{R}^{-1}\mathbf{C} \\ \hat{\mathbf{B}}^T \hat{\mathbf{B}} &= \mathbf{I} & \hat{\mathbf{C}} \hat{\mathbf{C}}^T &= \Lambda\end{aligned}$$

- R is computed as a generalized  $k \times k$  eigenvalue problem.  
(de la Torre, 2006)

$$(\mathbf{CC}^T)^{-1} \mathbf{R} = \mathbf{B}^T \mathbf{B} \mathbf{R} \Lambda^{-1}$$

## PCA/SVD in computer vision

- PCA/SVD has been applied to:
  - Recognition (eigenfaces: Turk & Pentland, 1991; Sirovich & Kirby, 1987; Leonardis & Bischof, 2000; Gong et al., 2000; McKenna et al., 1997a)
  - Parameterized motion models (Yacoob & Black, 1999; Black et al., 2000; Black, 1999; Black & Jepson, 1998)
  - Appearance/shape models (Cootes & Taylor, 2001; Cootes et al., 1998; Pentland et al., 1994; Jones & Poggio, 1998; Casia & Sclaroff, 1999; Black & Jepson, 1998; Blanz & Vetter, 1999; Cootes et al., 1995; McKenna et al., 1997; de la Torre et al., 1998b; de la Torre et al., 1998b)
  - Dynamic appearance models (Soatto et al., 2001; Rao, 1997; Oriols & Biñefica, 2001; Gong et al., 2000)
  - Structure from Motion (Tomasi & Kanade, 1992; Bregler et al., 2000; Sturm & Triggs, 1996; Brand, 2001)
  - Illumination based reconstruction (Hayakawa, 1994)
  - Visual servoing (Murase & Nayar, 1995; Murase & Nayar, 1994)
  - Visual correspondence (Zhang et al., 1995; Jones & Malik, 1992)
  - Camera motion estimation (Hartley, 1992; Hartley & Zisserman, 2000)

## More PCA/SVD work

- PCA/SVD has been applied to:
  - Image watermarking (Liu & Tan, 2000)
  - Signal processing (Moonen & de Moor, 1995)
  - Neural approaches (Oja, 1982; Sanger, 1989; Xu, 1993)
  - Bilinear models (Tenenbaum & Freeman, 2000; Marimont & Wandell, 1992)
  - Direct extensions (Welling et al., 2003; Penev & Atick, 1996)
- And many more (google)...
  - Results 1 - 10 of about 1,870,000 for "principal component analysis".
- Still work to do
  - Results 1 - 10 of about 65,300,000 for "Britney spears".

## 1-Robust PCA

- Two types of outliers:



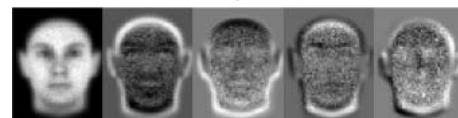
Sample outliers

(Xu & Yuille, 1995)

Intra-sample outliers

(de la Torre & Black, 2001b; Skocaj & Leonardis, 2003)

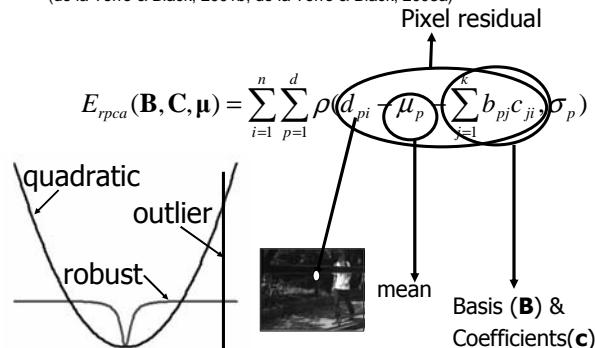
- Standard PCA solution (noisy data):



## Robust PCA

- Using robust statistics:

(de la Torre & Black, 2001b; de la Torre & Black, 2003a)



## Numerical problems

- No closed form solution in terms of an eigen-equation.
- Deflation approaches do not hold.

$$\mathbf{A}' = \mathbf{A} - \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T \quad \text{First eigenvector with highest eigenvalue.}$$

$$\mathbf{A}'' = \mathbf{A}' - \lambda_2 \mathbf{u}_2 \mathbf{u}_2^T \quad \text{Second eigenvector with highest eigenvalue.}$$

...

- In the robust case all the basis have to be computed simultaneously (including the mean).

## How to optimize it?

$$E_{rpca}(\mathbf{B}, \mathbf{C}, \boldsymbol{\mu}) = \sum_{i=1}^n \sum_{p=1}^d \rho(d_{pi} - \mu_p - \sum_{j=1}^k b_{pj}c_{ji}, \sigma_p)$$

- Normalized Gradient descent

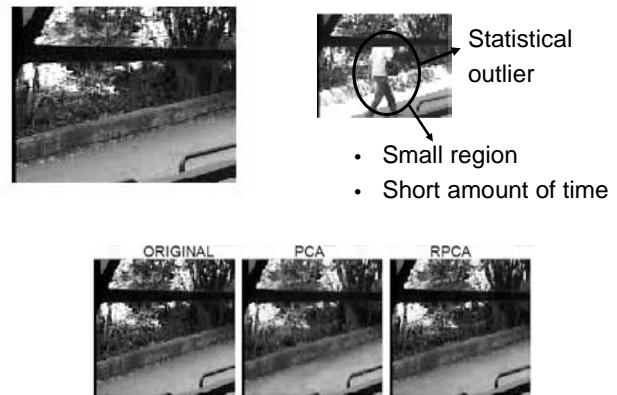
$$\mathbf{B}^{n+1} = \mathbf{B}^n - [\mathbf{H}_b]^{-1} \circ \frac{\partial E_{rpca}}{\partial \mathbf{B}} \quad \mathbf{H}_b = \max \operatorname{diag} \left( \frac{\partial^2 E_{rpca}}{\partial \mathbf{b}_i \partial \mathbf{b}_i^T} \right)$$

$$\mathbf{C}^{n+1} = \mathbf{C}^n - [\mathbf{H}_c]^{-1} \circ \frac{\partial E_{rpca}}{\partial \mathbf{C}} \quad \mathbf{H}_c = \max \operatorname{diag} \left( \frac{\partial^2 E_{rpca}}{\partial \mathbf{c}_i \partial \mathbf{c}_i^T} \right)$$

- Deterministic annealing methods to avoid local minima.

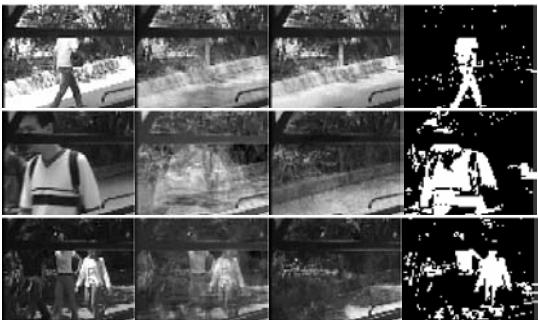
(Blake & Zisserman, 1987)

## Example

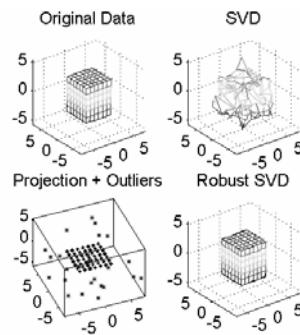


## Robust PCA

Original    PCA    RPCA    Outliers



## Structure from Motion



## Related RPCA work

- Robust estimation of coefficients  
(Black & Jepson, 1998; Leonardis & Bischof, 2000; Ke & Kanade, 2004)
- Robust estimation of basis and coefficients  
(Gabriel & Odoro, 1984; Croux & Filzmoser, 1981; Skocaj et al., 2002; Skocaj & Leonardis, 2003; de la Torre & Black, 2001b; de la Torre & Black, 2003a)
- Other Robust PCA techniques (sample outliers)  
(Campbell, 1980; Ruymagaart, 1981; Xu & Yuille., 1995)

## 2- PCA with uncertainty and missing data

- Adding uncertainty  $E_2(\mathbf{B}, \mathbf{C}) = \|\mathbf{W} \circ (\mathbf{D} - \mathbf{BC})\|_F^2 = \sum_{i=1}^d \sum_{j=1}^n w_j (d_{ij} - \sum_{s=1}^k b_{is} c_{sj})^2$



- If weights are separable  $\mathbf{W} = \mathbf{w}_r \mathbf{w}_c^T$  close form solution.

$$\mathbf{w}^c = \begin{pmatrix} w_1^c & w_2^c & \dots & w_n^c \end{pmatrix}$$

$$\mathbf{w}^r = \begin{pmatrix} w_1^r \\ w_2^r \\ \vdots \\ w_d^r \end{pmatrix}$$

$$\mathbf{D} = \begin{pmatrix} d_{11} & \dots & d_{1n} \\ d_{21} & \ddots & d_{2n} \\ \vdots & \ddots & \vdots \\ d_{d1} & \dots & d_{dn} \end{pmatrix}$$

$\mathbf{W} \in \Re^{d \times n}$        $w_{ij} \geq 0$

◦ Hadamard product

- Generalized SVD  
(Greenacre, 1984; Irani & Anandan, 2000;)

## General case

- For arbitrary weights no closed-form solution.  

$$E_2(\mathbf{B}, \mathbf{C}) = \|\mathbf{W} \circ (\mathbf{D} - \mathbf{BC})\|_F = \sum_{i=1}^d (\mathbf{d}_i - \mathbf{B}\mathbf{c}_i)^T \text{diag}(\mathbf{w}_i) (\mathbf{d}_i - \mathbf{B}\mathbf{c}_i) =$$
  

$$\sum_{i=1}^d (\mathbf{d}^p - \mathbf{C}^T \mathbf{b}^p)^T \text{diag}(\mathbf{w}^p) (\mathbf{d}^p - \mathbf{C}^T \mathbf{b}^p)$$
 (Torre & Black, 2003a)

- Alternated least squares algorithms

- Slow convergence, easy implementation.

- Damped Newton Algorithm

- Fast convergence. (Buchanan & Fitzgibbon., 2005)

$$E_2(\mathbf{B}, \mathbf{C}) = \|\mathbf{W} \circ (\mathbf{D} - \mathbf{BC})\|_F + \lambda_1 \|\mathbf{B}\|_F + \lambda_2 \|\mathbf{C}\|_F$$

$$\mathbf{v} = \begin{bmatrix} \text{vec}(\mathbf{B}) \\ \text{vec}(\mathbf{C}) \end{bmatrix} \quad \mathbf{v}^{(n+1)} = \mathbf{v}^n - \left[ \frac{\partial^2 E_2}{\partial^2 \mathbf{v}} \right]^{-1} \frac{\partial E_2}{\partial \mathbf{v}}$$

$$-\mathbf{H} \text{ definite positive: } \mathbf{H} = \frac{\partial^2 E_2}{\partial^2 \mathbf{v}} + \lambda \mathbf{I}$$

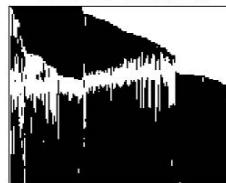
```

repeat
   $\mathbf{H} = \frac{\partial^2 E_2}{\partial^2 \mathbf{v}} \quad \mathbf{g} = \frac{\partial E_2}{\partial \mathbf{v}}$ 
repeat
   $\lambda = 10\lambda$ 
   $\mathbf{y} = \mathbf{x} - (\mathbf{H} + \lambda I)^{-1} \mathbf{g}$ 
until  $F(\mathbf{y}) < F(\mathbf{x})$ 
 $\mathbf{x} = \mathbf{y}; \lambda = \frac{\lambda}{10}$ 
until convergence
  
```

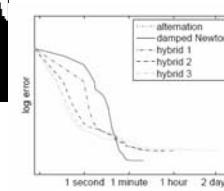
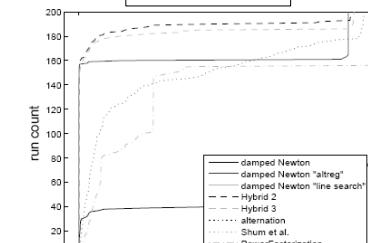
## Experiments



240 × 167: 70% known



Total: 500 runs



## Related work

- Iterative (Wiberg, 1976; Shum et al., 1995; Morris & Kanade, 1998; Aans et al., 2002; Guerreiro & Aguilar, 2002)
- Closed-form (Aguiar & Moura, 1999; Irani & Anandan, 2000)
- Power factorization (Hartley & Schaalitzky, 2003)
- Bayesian estimation (Torresani & Bregler, 2004)

### Incremental PCA

- (de la Torre et al., 1998b; Ross et al., 2004; Brand, 2002; Skocaj & Leonardis, 2003; Champagne & Liu, 1998; A. Levy, 2000)

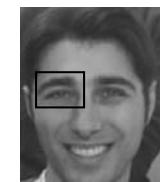
## 3- Parameterized Component Analysis (PaCA) (de la Torre & Black, 2003b)

- Learn a subspace invariant to geometric transformations?

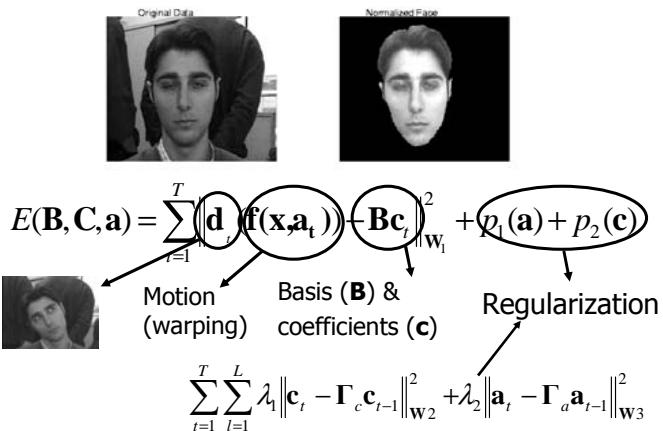


- Data has to be **geometrically** normalized

- Tedious manual cropping.
- Inaccuracies due to matching ambiguities.
- Hard to achieve sub-pixel accuracy.



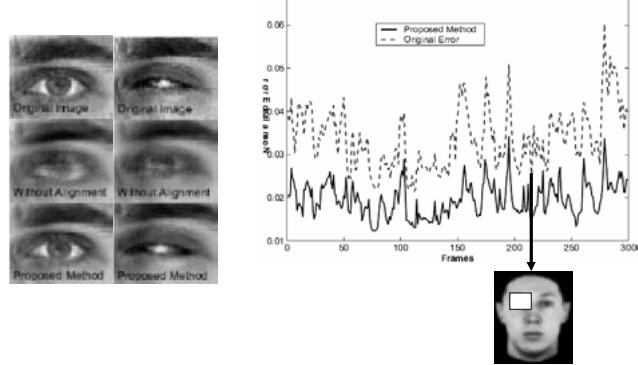
## Error function for PaCA



## Solving the optimization problem

- Linearizing the motion (Bergen et al., 1992; Black & Jepson, 1998)
 
$$\mathbf{d}_t(\mathbf{f}(\mathbf{x}, \mathbf{a}_t)) \approx \mathbf{d}_t(\mathbf{f}(\mathbf{x}, \mathbf{a}_{0t})) + \mathbf{J}_t \Delta \mathbf{a}_t$$
- Normalized gradient descent w.r.t. all parameters + deterministic annealing.
  - Update for  $\mathbf{c}$  (appearance) &  $\mathbf{a}$  (motion).
  - Updated for  $\mathbf{B}$  (appearance basis).
- It is a non-convex function.
  - Stochastic initialization (G.A). (Lanitis et al., 1995; de la Torre & Black, 2003b)
  - Multiresolution motion estimation framework.

## EigenEye Learning



## More on parameterized CA

- Probabilistic model
  - Search scales exponentially with the number of motion parameters (Frey & Jovic, 1999a; Frey & Jovic, 1999b; Williams & Titsias, 2004)
- Other continuous approaches.
  - (Schewitzer, 1999; Rao, 1999; Shashua et al., 2002)
- Invariant clustering
  - (Fitzgibbon & Zisserman, 2003)
- Non-rigid motion
  - (Baker et al., 2004)
- Invariant recognition
  - (Black & Jepson, 1998)
- Invariant support vector machines (Avidan, 2001)

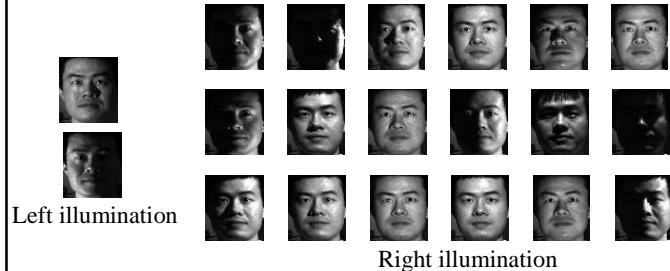
## 4- PCA over continuous spaces

(Levin & Shashua, 2002)

- PCA assumes discrete samples, but in the limit:

$$\frac{1}{N} \sum_{i=1}^N \mathbf{d}_i \mathbf{d}_i^T \xrightarrow[N \rightarrow \infty]{} \int f(\mathbf{d}) \mathbf{d} \mathbf{d}^T d\mathbf{d}$$

- Sometimes not uniform sampling of  $f(\mathbf{d})$ .



Component Analysis for Computer Vision

F. De la Torre

ECCV-06

29

## Bias solution

- 3 principal components (38 people)



Original

Bias

Unbiased



Component Analysis for Computer Vision

F. De la Torre

ECCV-06

30

## How to unbiased the solution?

- Weighting the data.

– Not clear which is the optimal weight.  $E(\mathbf{B}, \mathbf{C}) = \sum_{i=1}^n w_i \|\mathbf{d}_i - \mathbf{B}\mathbf{c}_i\|_F$

- More elegant approach.

– Find the principal components over the data points represented by (dense) uniform sampling.

– Integrating over the convex combination of the examples is less sensitive to the particular sample



Component Analysis for Computer Vision

F. De la Torre

ECCV-06

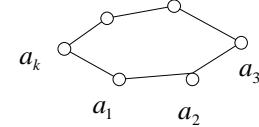
31

## PCA OVER POLYTOPS

- PCA over polytops

$$Cov(W) = \frac{1}{V(W)} \int_W aa^T da$$

Volume of polytopes



- A simple 2D case, two points  $\mathbf{a}_1$  and  $\mathbf{a}_2$ :

– Regular PCA  $\mathbf{a}_1 \mathbf{a}_1^T + \mathbf{a}_2 \mathbf{a}_2^T$   
– Integrating over a line  $\lambda \mathbf{a}_1 + (1 - \lambda) \mathbf{a}_2$ ,  $0 \leq \lambda \leq 1$

$$\max_{\|\mathbf{u}\|=1} \int_{\mathbf{a} \in W} |\mathbf{a}^\top \mathbf{u}|^2 d\mathbf{a} = \mathbf{u}^\top \left[ \int_{\mathbf{a} \in W} \mathbf{a} \mathbf{a}^\top d\mathbf{a} \right] \mathbf{u} \quad \int_0^1 \lambda^2 d\lambda = \frac{1}{3}, \quad \int_0^1 \lambda(1-\lambda) d\lambda = \frac{1}{6},$$

$$\max_{\|\mathbf{u}\|=1} \mathbf{u}^\top [\mathbf{a}_1 \mathbf{a}_1^\top + \mathbf{a}_2 \mathbf{a}_2^\top + \frac{1}{2} (\mathbf{a}_1 \mathbf{a}_2^\top + \mathbf{a}_2 \mathbf{a}_1^\top)] \mathbf{u} \quad A = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix} A^\top = A \Phi A^\top$$

Component Analysis for Computer Vision

F. De la Torre

ECCV-06

32

## In general

$$\Phi_k = \frac{1}{k(k+1)} \begin{bmatrix} 2 & 1 & \dots & 1 \\ 1 & 2 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 2 \end{bmatrix} = \frac{1}{k(k+1)} (\mathbf{I} + ee^T)$$

$$Cov(W) = \frac{1}{V(W)} \int_{a \in W} \mathbf{d} \mathbf{d}^T d\mathbf{d} = \frac{1}{k(k+1)} \mathbf{D} (\mathbf{I} + ee^T) \mathbf{D}^T$$

- Traditional PCA  $\frac{1}{k} \mathbf{D} \mathbf{D}^T$
- Continuous PCA  $\frac{1}{k(k+1)} \mathbf{D} (\mathbf{I} + ee^T) \mathbf{D}^T$

## 5-Filtered PCA

(Bischof et al., 2004; Wildenauer et al., 2002)

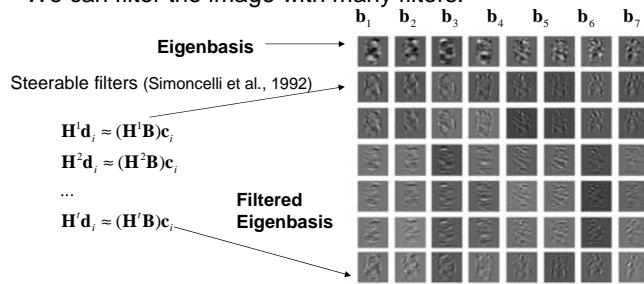
- How to construct eigenspaces robust:
  - Varying illumination.
  - Occlusion.
  - Noise.
- Filtered PCA.
  - $\mathbf{H}$  is a convolution matrix (block circulant structure)
  - The coefficients of  $\mathbf{c}_i$  remain the same under a convolution.



$$\mathbf{H} \in \Re^{d \times d} \quad \mathbf{d}_i \approx \mathbf{B} \mathbf{c}_i \Rightarrow \mathbf{H} \mathbf{d}_i \approx (\mathbf{H} \mathbf{B}) \mathbf{c}_i$$

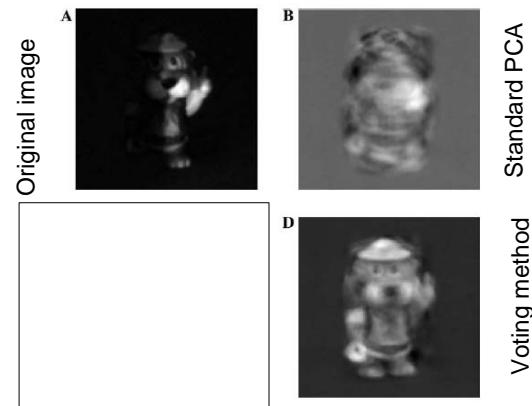
## Filtered PCA

- We can filter the image with many filters:



- Given a new image, compute the filtered representation and robustly compute  $\mathbf{c}_i$  (eigenbasis known).

## Experiments of Filtered PCA



## More results

Obj.	1	2	3	4	5	6	7	8	9	10	%	Ang.
<i>Robust filtered method—all eigenvectors used</i>												
1	360	0	0	0	0	0	0	0	0	0	100.0	4.92
2	0	321	3	0	0	0	0	0	0	0	99.1	9.03
3	0	1	503	0	0	0	0	0	0	0	99.8	0.99
4	3	1	0	355	1	0	0	0	0	0	98.6	3.30
5	0	0	0	0	612	0	0	0	0	0	100.0	3.79
6	0	9	0	0	11	672	13	1	8	6	93.3	15.85
7	0	14	0	8	22	10	413	26	7	4	81.9	4.46
8	0	11	0	0	51	37	1	284	42	6	65.7	15.74
9	0	16	2	9	17	13	0	7	439	1	87.1	14.10
10	1	0	3	0	9	46	0	5	39	509	83.2	13.97
Avg.											90.6	8.71
<i>Standard method—all eigenvectors used</i>												
1	74	0	0	40	179	41	4	2	4	16	20.6	2.84
2	0	215	0	6	2	67	0	1	20	13	66.4	26.37
3	0	7	237	3	187	28	7	0	0	35	47.0	3.25
4	3	2	0	236	43	28	0	0	44	4	65.6	12.03
5	0	0	8	0	535	26	7	17	19	0	87.4	8.73
6	3	78	0	0	65	553	0	3	4	14	76.8	41.14
7	0	19	0	20	96	74	181	44	22	48	35.9	6.02
8	0	30	0	0	134	124	3	119	22	0	27.5	16.47
9	0	10	0	54	62	102	4	11	255	6	50.6	22.24
10	4	10	2	24	85	175	5	4	36	267	43.6	14.08
Avg.											54.2	19.48

Component Analysis for Computer Vision

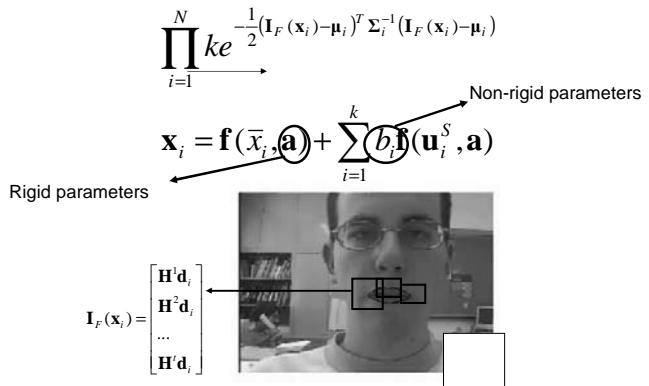
F. De la Torre

ECCV-06

37

## Eigenfiltering for flexible Eigentracking

(de la Torre et al., 2000)



Component Analysis for Computer Vision

F. De la Torre

ECCV-06

38

## 6- PCA of a set of Rotated Images

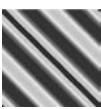
(Uenohara & Kanade, 1998; Jogan et al., 2003)

- Set of uniformly in-plane rotated versions of the same object (0 background).



- $\mathbf{F} = \mathbf{D}^T \mathbf{D}$  is a circulant symmetric Toeplitz matrix.

$$\mathbf{D}^T \mathbf{D} = \begin{bmatrix} \mathbf{d}_1^T \mathbf{d}_1 & \mathbf{d}_1^T \mathbf{d}_2 & \dots & \mathbf{d}_1^T \mathbf{d}_n \\ \mathbf{d}_2^T \mathbf{d}_1 & \mathbf{d}_2^T \mathbf{d}_2 & \dots & \mathbf{d}_2^T \mathbf{d}_n \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{d}_n^T \mathbf{d}_1 & \mathbf{d}_n^T \mathbf{d}_2 & \dots & \mathbf{d}_n^T \mathbf{d}_n \end{bmatrix} = \begin{bmatrix} d_0 & d_1 & \dots & d_{n-1} \\ d_{n-1} & d_0 & d_1 & d_{n-2} \\ \vdots & \vdots & \ddots & \vdots \\ d_{n-2} & d_{n-1} & \dots & \dots \end{bmatrix}$$



Component Analysis for Computer Vision

F. De la Torre

ECCV-06

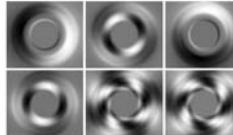
39

## Eigenvectors of circulant matrices

- The eigenvectors of circulant matrices are the n basis of the Fourier matrix.

$$\mathbf{F} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & e^{-\frac{2\pi i}{n}} & \dots & e^{-\frac{(n-1)*2\pi i}{n}} \\ 1 & e^{-\frac{2\pi i}{n}} & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots \\ 1 & e^{-\frac{(n-2)*2\pi i}{n}} & \dots & e^{-\frac{(n-1)*2\pi i}{n}} \end{bmatrix} \quad \mathbf{V} = \begin{bmatrix} \sum_{k=0}^{n-1} d_k e^{-\frac{k*2\pi i}{n}} \\ \sum_{k=0}^{n-1} d_k e^{-\frac{(k+1)*2\pi i}{n}} \\ \vdots \\ \sum_{k=0}^{n-1} d_k e^{-\frac{(k+n-1)*2\pi i}{n}} \end{bmatrix} \quad i = \sqrt{-1} \quad \mathbf{F}\mathbf{F}^H = \mathbf{V}\mathbf{V}^H$$

- Circulant and SYMMETRIC  $d_i = d_{n-i}$ .  $\mathbf{F}$  is the basis of cosines.



eig(DD^T) → DCT

Component Analysis for Computer Vision

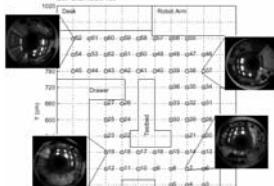
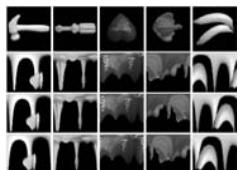
F. De la Torre

ECCV-06

40

## Generalization to multiple templates

(Jogan et al., 2003)



- $P$  different locations (objects), each shifted  $N$  times
- every  $Q_{ij}$  is circulant (but in general not symmetric!)

$$\mathbf{A} = \mathbf{D}' \mathbf{D} = \begin{bmatrix} Q_{00} & \dots & Q_{0(P-1)} \\ Q_{10} & \dots & \dots \\ \dots & \dots & \dots \\ Q_{(P-1)0} & \dots & Q_{(P-1)(P-1)} \end{bmatrix}$$



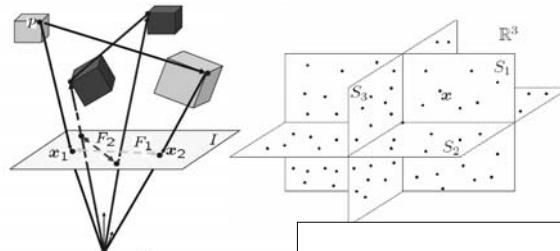
Eigencalculation of  $O(NP^3)$  rather than  $O((NP)^3)$

## 7- Mixture of subspaces

(Vidal et al., 2003)

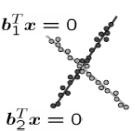
- How to estimate (basis) of a mixture of subspaces and varying dimensions?

- Multibody structure from motion
  - 2D & 3D Motion segmentation.



## Generalized Principal Component Analysis

- Given points on multiple subspaces, identify:
  - Number of subspaces and their dimensions
  - Basis of each subspace
  - Segmentation of data points.
- Classical chicken-egg problem.
- GPCA finds an algebraic geometric solution to the mixture of subspaces.
  - Number of subspaces= degrees of polynomial
  - Subspaces basis= derivatives of polynomial.



## Example: clustering in 1D

$$x = b_1 \quad x = b_2$$

$$x = b_1 \text{ or } x = b_2$$

$$(x - b_1)(x - b_2) = 0$$

$$x^2 - (b_1 + b_2)x + b_1 b_2 = 0$$

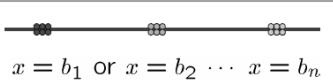
$$\begin{bmatrix} x_1^2 & x_1 & 1 \\ x_2^2 & x_2 & 1 \\ \vdots & \vdots & \vdots \\ x_N^2 & x_N & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -(b_1 + b_2) \\ b_1 b_2 \end{bmatrix} = 0$$

- Number of groups?

$\text{rank}(P) = 1$ : one group only

$\text{rank}(P) = 2$ : two groups

## Example 2



$$p_n(x) = (x - b_1) \cdots (x - b_n) = 0$$

$$p_n(x) = x^n + c_1 x^{n-1} + \cdots + c_n = 0$$

$$p_n(x) = [x^n \quad \cdots \quad x \quad 1] c = 0$$

$$P_n c = \begin{bmatrix} x_1^n & \cdots & x_1 & 1 \\ x_2^n & \cdots & x_2 & 1 \\ \vdots & & \vdots & \vdots \\ x_N^n & \cdots & x_N & 1 \end{bmatrix} c = 0$$

$P_n \in \mathbb{R}^{N \times (n+1)}$

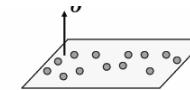
How to compute n, c, b's?

- Number of clusters
- $n = \min\{i : \text{rank}(P_i) = i\}$
- Cluster centers  
Roots of  $p_n(x)$
- Solution is unique if  
 $N_{\text{points}} \geq n_{\text{groups}}$
- Solution is closed form if  
 $n_{\text{groups}} \leq 4$

## GPCA

- One plane

$$b^T x = b_1 x_1 + b_2 x_2 + b_3 x_3 = 0$$



- One line

$$\begin{aligned} b_1^T x &= b_1 x_1 + b_2 x_2 + b_3 x_3 = 0 \\ b_2^T x &= b_4 x_1 + b_5 x_2 + b_6 x_3 = 0 \end{aligned}$$

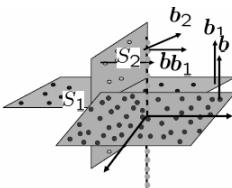
- One subspace can be represented with

- Set of linear equations  $S = \{\mathbf{x} : B^T \mathbf{x} = 0\}$
- Set of polynomials of degree 1

## GPCA

- Two planes  
 $(b_1^T x = 0) \text{ or } (b_2^T x = 0)$

$$p_2(x) = (b_1^T x)(b_2^T x) = 0$$



- One plane and one line

- Plane:  $S_1 = \{\mathbf{x} : b_1^T \mathbf{x} = 0\}$
- Line:  $S_2 = \{\mathbf{x} : b_1^T \mathbf{x} = b_2^T \mathbf{x} = 0\}$

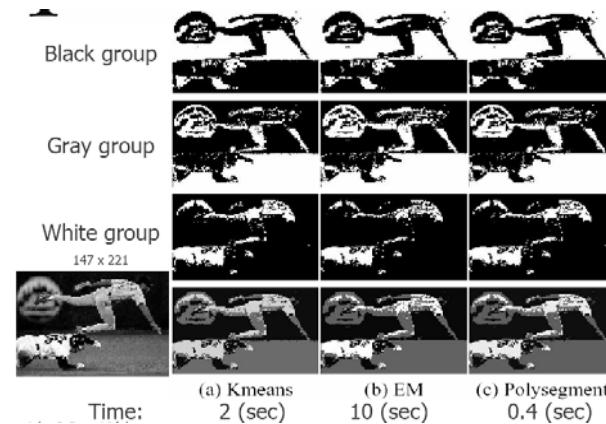
$$S_1 \cup S_2 = \{\mathbf{x} : (b_1^T \mathbf{x}) = 0 \text{ or } (b_1^T \mathbf{x}) = b_2^T \mathbf{x} = 0\}$$

De Morgan's rule

$$S_1 \cup S_2 = \{\mathbf{x} : (b_1^T \mathbf{x})(b_1^T \mathbf{x}) = 0 \text{ and } (b_1^T \mathbf{x})(b_2^T \mathbf{x}) = 0\}$$

- A union of n subspaces can be represented with a set of homogeneous polynomials of degree n

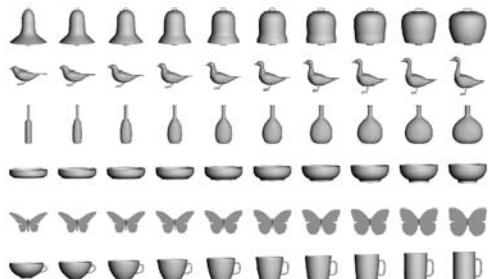
## Examples



## Multiple eigenspaces

(Leonardis et al., 2002)

- Goal: group visually similar images into categories in an unsupervised and self-organising way.



## Multiple eigenspaces – idea

- Group the images and construct multiple eigenspaces such that:
  - Mean reconstruction error is always below a threshold
  - The dimensionalities of eigenspaces are as small as possible
- Start from a seed and then grow hypotheses
- Grow and select paradigm:
  - Simultaneously and independently grow multiple competing hypotheses
  - Select the best hypotheses at the end



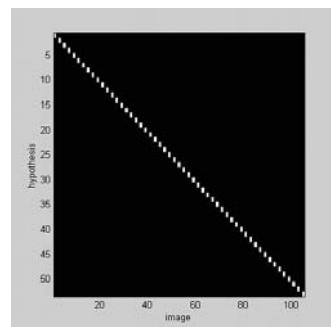
## Multiple eigenspace - example

- 5 categories, 21 images from each



Growing

:



## Multiple eigenspace - example

- Results:

mean	eigenvectors
------	--------------

corresponding images

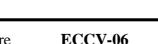
ES1:     

ES2:     

ES3:    

ES4:   

ES5:  



## "Intercorrelations among variables are the bane of the multivariate researcher's struggle for meaning"

Cooley and Lohnes, 1971



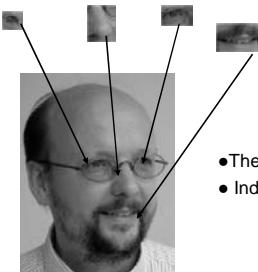
Component Analysis for Computer Vision

F. De la Torre

ECCV-06

53

## Part-based representation



- The firing rates of neurons are never negative
- Independent representations.

NMF & ICA

Component Analysis for Computer Vision

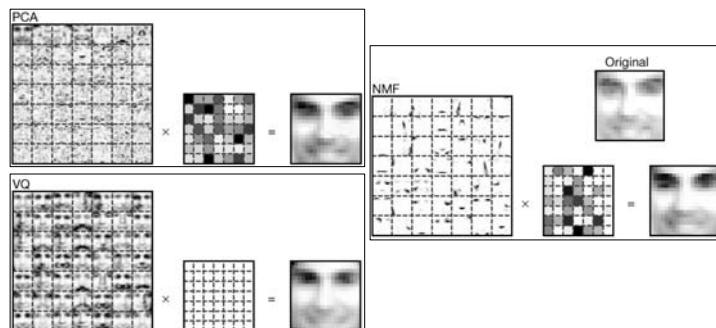
F. De la Torre

ECCV-06

54

## Non-negative Matrix Factorization

- Positive factorization.
- $E(\mathbf{B}, \mathbf{C}) = \|\mathbf{D} - \mathbf{BC}\|_F$     $\mathbf{B}, \mathbf{C} \geq 0$
- Leads to part-based representation.



Component Analysis for Computer Vision

F. De la Torre

ECCV-06

55

## Nonnegative factorization algorithm

(Lee & Seung, 1999; Lee & Seung, 2000)

$$\min_{\mathbf{W} \geq 0, \mathbf{V} \geq 0} F = \sum_{ij} |d_{ij} - (\mathbf{BC})_{ij}|^2$$

Inference:

$$\mathbf{C}_{ij} \leftarrow \mathbf{C}_{ij} \frac{(\mathbf{B}^T \mathbf{D})_{ij}}{(\mathbf{B}^T \mathbf{B} \mathbf{V})_{ij}}$$

Derivatives:

$$\frac{\partial F}{\partial \mathbf{C}_{ij}} = (\mathbf{B}^T \mathbf{B} \mathbf{C})_{ij} - (\mathbf{B}^T \mathbf{C})_{ij}$$

Learning:

$$\frac{\partial F}{\partial \mathbf{B}_{ij}} = (\mathbf{B} \mathbf{C}^T)_{ij} - (\mathbf{D} \mathbf{C}^T)_{ij}$$

- Multiplicative algorithm can be interpreted as diagonally rescaled gradient descent.

Component Analysis for Computer Vision

F. De la Torre

ECCV-06

## 8- PCA & NMF for clustering

(Zha et al., 2001; Ding & He, 2004; de la Torre & Kanade, 2006)

- VQ, k-means  $E(\mathbf{G}, \mathbf{B}) = \|\mathbf{D} - \mathbf{BG}^T\|_F = \sum_{i=1}^c \sum_{j \in C_i} \|\mathbf{d}_j - \mathbf{b}_i\|$
- $$\mathbf{G}^T = \begin{bmatrix} 1 & \dots & 0 \\ 0 & \dots & 1 \\ 0 & \dots & 0 \end{bmatrix} \quad g_{ij} \in \{0,1\} \quad \mathbf{G}\mathbf{1}_k = \mathbf{1}_n \quad \mathbf{G} \in \mathbb{R}^{n \times k}$$

- After eliminating  $(\mathbf{B})$

$$E(\mathbf{G}) = \|\mathbf{D} - \mathbf{DG}(\mathbf{G}^T \mathbf{G})^{-1} \mathbf{G}^T\|_F = \text{tr}(\mathbf{D}^T \mathbf{D}) - \text{tr}((\mathbf{G}^T \mathbf{G})^{-1} \mathbf{G}^T \mathbf{D}^T \mathbf{D} \mathbf{G}) \geq \sum_{i=c+1}^{\min(d,n)} \lambda_i$$

- Relaxing binary constraints

$$E(\mathbf{G}) \propto \text{tr}((\mathbf{G}^T \mathbf{G})^{-1} (\mathbf{G}^T \mathbf{D}^T \mathbf{D} \mathbf{G}))$$

- Eigenvectors (PCA) of  $\mathbf{D}^T \mathbf{D}$  are the optimal continuous solution of the indicator variable  $\mathbf{G}$ .

## Clustering with NMF

(Zass & Shashua, 2005; Ding et al., 2005)

$$E(\mathbf{G}, \mathbf{B}) = \|\mathbf{D} - \mathbf{BG}^T\|_F \quad \mathbf{G} = [\varphi(\mathbf{d}_1) \ \varphi(\mathbf{d}_2) \ \dots \ \varphi(\mathbf{d}_n)]$$

- Soft clustering and non-negative matrix factorization:

$$\left\| \mathbf{\Gamma}^T \mathbf{\Gamma} - \mathbf{GG}^T \right\|_F \quad \mathbf{G} \geq 0 \quad \mathbf{F} = \mathbf{G}^T \mathbf{G} = \mathbf{I} \quad \mathbf{F} \mathbf{1} = \mathbf{1} \quad \mathbf{F}^T \mathbf{1} = \mathbf{1}$$

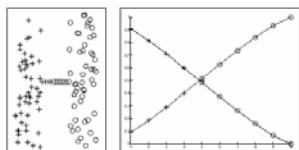
Affinity matrix

$$\mathbf{G} = \begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{\sqrt{n_1}} & \frac{1}{\sqrt{n_2}} & \dots & \frac{1}{\sqrt{n_c}} \\ \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

- Previous normalization  $\min_{\mathbf{K}} \|\mathbf{\Gamma}^T \mathbf{\Gamma} - \mathbf{K}\|_F$  s.t.  $\mathbf{K} \mathbf{1} = \mathbf{K}^T \mathbf{1} = \beta \mathbf{1}$   $\mathbf{K} = \mathbf{K}^T$
- Gradient descent with multiplicative updates.

## Examples

- Soft clustering



$R_{ij}=1$  if  $\mathbf{d}_i$  and  $\mathbf{d}_j$  are cluster together.  
 $R_{ij}=-1$  if  $\mathbf{d}_i$  and  $\mathbf{d}_j$  are not cluster together.

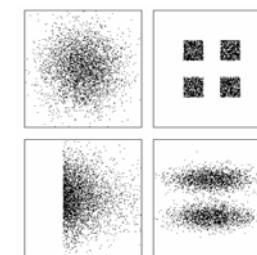
- Adding side-information

$$\mathbf{K}' = \mathbf{K} + \alpha \mathbf{R} \quad (\text{Zass & Shashua, 2005})$$

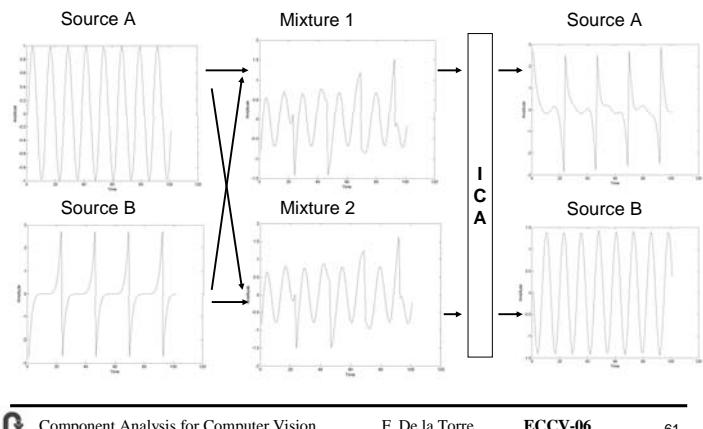


## Independent Component Analysis

- We need more than second order statistics to represent the signal.



## ICA & Signals



Component Analysis for Computer Vision

F. De la Torre

ECCV-06

61

## ICA

(Hyvriinen et al., 2001)

$$\mathbf{D} = \mathbf{B}\mathbf{C} \quad \mathbf{C} \approx \mathbf{S} = \mathbf{W}\mathbf{D} \quad \mathbf{W} \approx \mathbf{B}^{-1}$$

- Look for  $s_i$  that are independent.
- PCA finds uncorrelated variables, the independent components have non Gaussian distributions.
- Uncorrelated  $E(s_i s_j) = E(s_i)E(s_j)$
- Independent  $E(g(s_i)f(s_j)) = E(g(s_i))E(f(s_j))$  for any non-linear  $f,g$



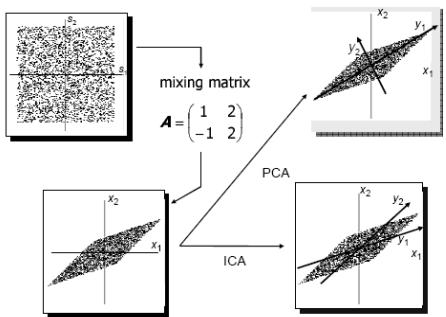
Component Analysis for Computer Vision

F. De la Torre

ECCV-06

62

## ICA vs PCA



Component Analysis for Computer Vision

F. De la Torre

ECCV-06

63

## Many optimization criteria

- Minimize high order moments: e.g. kurtosis  

$$\text{kurt}(\mathbf{W}) = E\{s^4\} - 3(E\{s^2\})^2$$
- Many other information criteria.
- Also an error function: (Olhausen & Field, 1996)

$$\sum_{i=1}^n \|\mathbf{d}_i - \mathbf{B}\mathbf{c}_i\| + \sum_{i=1}^n S(\mathbf{c}_i) \quad \text{Sparseness (e.g. } S=|\cdot|)$$

- Other sparse PCA.

(Chennubhotla & Jepson, 2001b; Zou et al., 2005; dAspremont et al., 2004;)

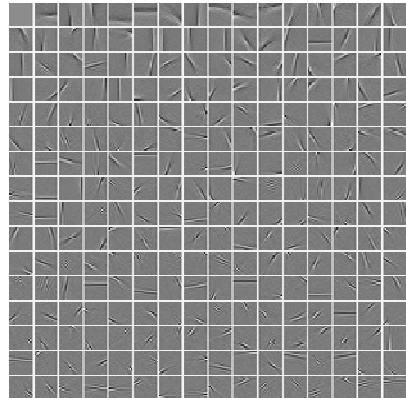
Component Analysis for Computer Vision

F. De la Torre

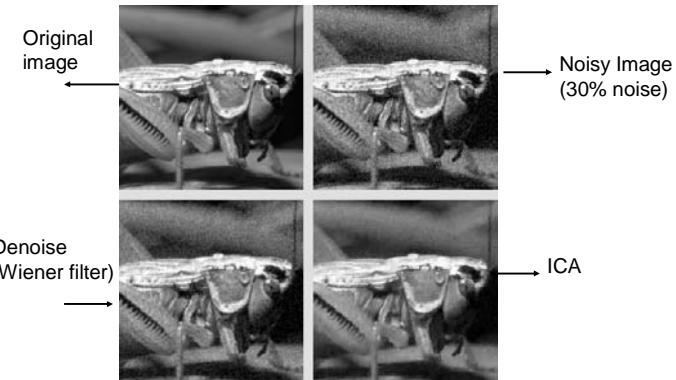
ECCV-06

64

## Basis of natural images



## Denoising



## Discriminative Models

- Linear Discriminant Analysis (LDA)
  - 9) Multimodal Oriented Discriminant Analysis.
  - 10) Discriminative Cluster Analysis.
  - 11) Robust Linear Discriminant Analysis.
- Oriented Component Analysis (OCA)
  - 12) Representational Oriented Component Analysis.
- Canonical Correlation Analysis (CCA)
  - 13) Dynamical Coupled Component Analysis.
  - 14) CCA and Mobile robotics applications.
- Relevance Component Analysis (RCA)

## Linear Discriminant Analysis (LDA)

(Fisher, 1938; Mardia et al., 1979; Bishop, 1995)

$$S_b = \sum_{i=1}^c \sum_{j=1}^c (\mu_i - \mu_j)(\mu_i - \mu_j)^T$$

$$S_w = \sum_{j=1}^{C_i} \sum_{i=1}^{C_j} (\mathbf{d}_i - \mu_j)(\mathbf{d}_i - \mu_j)^T$$

$$S_t = \mathbf{D}\mathbf{D}^T = \sum_{i=1}^n \mathbf{d}_i \mathbf{d}_i^T$$

$$J(\mathbf{B}) = \frac{|\mathbf{B}^T S_b \mathbf{B}|}{|\mathbf{B}^T S_w \mathbf{B}|}$$

$$\mathbf{S}_b \mathbf{B} = \mathbf{S}_t \mathbf{B} \Lambda$$

- Optimal linear dimensionality reduction if classes are Gaussian with equal covariance matrix.

## Matrix formulation

(de la Torre & Kanade, 2005b)

$$\mathbf{S}_t = \frac{1}{n-1} \sum_{j=1}^n (\mathbf{d}_j - \mathbf{m})(\mathbf{d}_j - \mathbf{m})^T = \frac{1}{n-1} \mathbf{D} \mathbf{P}_1 \mathbf{D}^T$$

$$\mathbf{S}_w = \frac{1}{n-1} \sum_{i=1}^c \sum_{d_j \in C_i} (\mathbf{d}_j - \mathbf{m}_i)(\mathbf{d}_j - \mathbf{m}_i)^T = \frac{1}{n-1} \mathbf{D} \mathbf{P}_2 \mathbf{D}^T$$

$$\mathbf{S}_b = \sum_{i=1}^c \frac{n_i}{n-1} (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T = \frac{1}{n-1} \mathbf{D} \mathbf{P}_3 \mathbf{D}^T$$

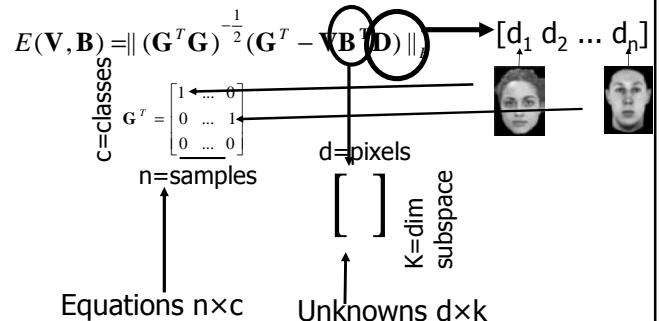
$$\mathbf{P}_1 = \mathbf{I} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \quad \mathbf{P}_2 = \mathbf{I} - \mathbf{G}(\mathbf{G}^T \mathbf{G})^{-1} \mathbf{G}^T$$

$$\mathbf{P}_3 = \mathbf{G}(\mathbf{G}^T \mathbf{G})^{-1} \mathbf{G}^T - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$$

- $\mathbf{m}$  overall mean,  $\mathbf{m}_i$  mean for class i.
- $\mathbf{P}_i$  is a projection matrix (i.e.  $\mathbf{P}_i = \mathbf{P}_i^T \quad \mathbf{P}_i = \mathbf{P}_i^2$ )

## Error function for LDA

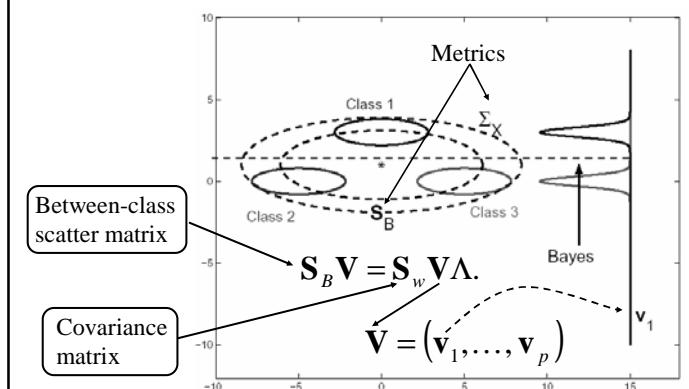
(de la Torre & Kanade, 2006)



- $d \gg n$  an UNDETERMINED system of equations! (overfitting)

## Where are linear methods applicable?

(Martinez & Zhu., 2005)



## When is LDA good?

(Martinez & Zhu., 2005)

- The *discriminant power* of the generalized eigenvalue decomposition equation

$$\mathbf{M}_A \mathbf{V} = \mathbf{M}_B \mathbf{V} \Lambda$$

is  $\text{tr}(\mathbf{M}_B^{-1} \mathbf{M}_A)$ , which is the same as

$$K = \sum_{i=1}^q \sum_{j=1}^p \frac{\lambda_{A_j}}{\lambda_{B_j}} (\mathbf{b}_j^T \mathbf{a}_i)^2.$$

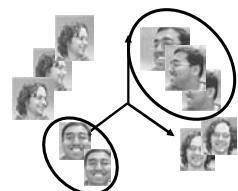
relative variance

- Large values of K indicates instability in the results.
- If for some (i,j)  $\max (\mathbf{b}_j^T \mathbf{a}_i)^2$  close to 1 the solution might be unstable.

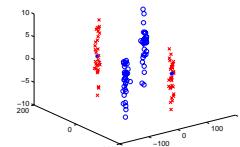
## 9- Multimodal Oriented Component Analysis (MODA)

(de la Torre & Kanade, 2005a)

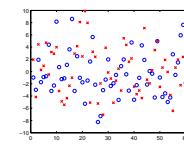
- How to extend LDA
  - Model class covariances.
  - Multimodal classes.
  - Deal efficiently with huge covariance matrices (e.g. 100\*100).



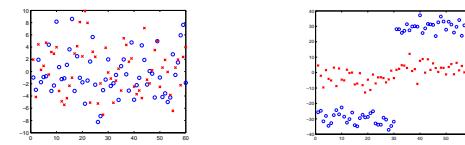
## Multimodality



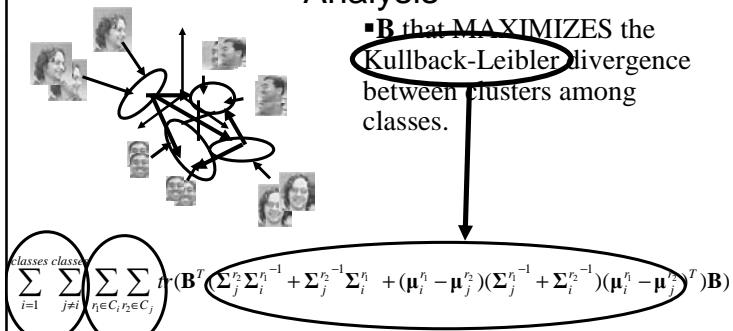
LDA



MODA



## Multimodal Oriented Discriminant Analysis



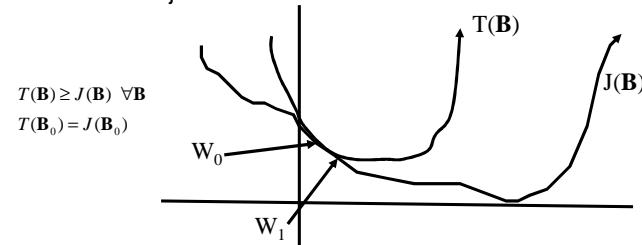
- 1 mode per class and equal covariances equivalent to LDA.

## Optimization

- Hard optimization problem

$$J(\mathbf{B}) = - \sum_{i=1}^{classes} \text{tr}((\mathbf{B}^T \Sigma_i \mathbf{B})^{-1} (\mathbf{B}^T \mathbf{A}_i \mathbf{B}))$$

- Iterative Majorization (Kiers, 1995; Leeuw, 1994)



## Majorization

$$T(\mathbf{B}) = \sum_{i=1}^{classes} \| (\mathbf{B}^T \Sigma_i \mathbf{B})^{-\frac{1}{2}} \mathbf{B}^T \mathbf{A}_i^{-\frac{1}{2}} - (\mathbf{B}^T \Sigma_i \mathbf{B})^{\frac{1}{2}} (\mathbf{B}_0^T \Sigma_i \mathbf{B}_0)^{-\frac{1}{2}} \mathbf{B}^T \mathbf{A}_i^{\frac{1}{2}} \|$$

$$\geq - \sum_{i=1}^{classes} \text{tr}((\mathbf{B}^T \Sigma_i \mathbf{B})^{-1} (\mathbf{B}^T \mathbf{A}_i \mathbf{B}))$$

- Slow convergence, first gradient descent:

$$\mathbf{B}^{(n+1)} = \mathbf{B}^{(n)} - \eta \frac{\partial E_7(\mathbf{B}^{(n)})}{\partial \mathbf{B}}$$

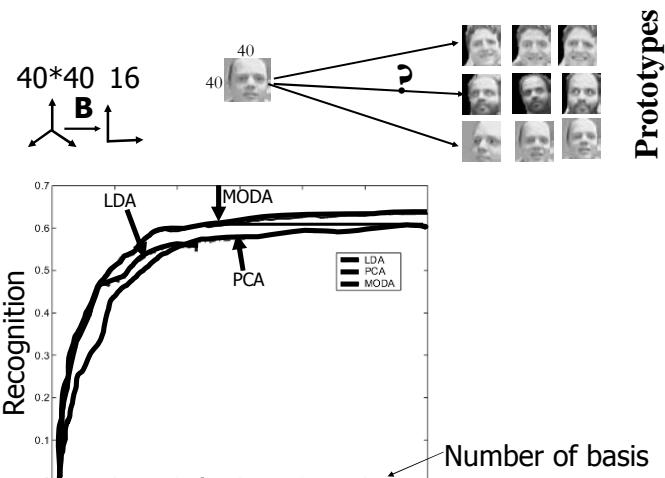
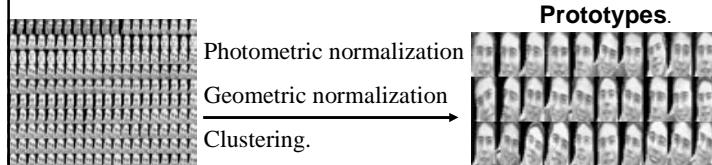
$$\frac{\partial E_7}{\partial \mathbf{B}} = \sum_i \mathbf{A}_i \mathbf{B}^{(n)} (\mathbf{B}^{(n)T} \Sigma_i \mathbf{B}^{(n)})^{-1} - \Sigma_i \mathbf{B}^{(n)} (\mathbf{B}^{(n)T} \Sigma_i \mathbf{B}^{(n)})^{-1} (\mathbf{B}^{(n)T} \mathbf{A}_i \mathbf{B}^{(n)})^{-1} (\mathbf{B}^{(n)T} \Sigma_i \mathbf{B}^{(n)})^{-1}$$

## Face recognition from video



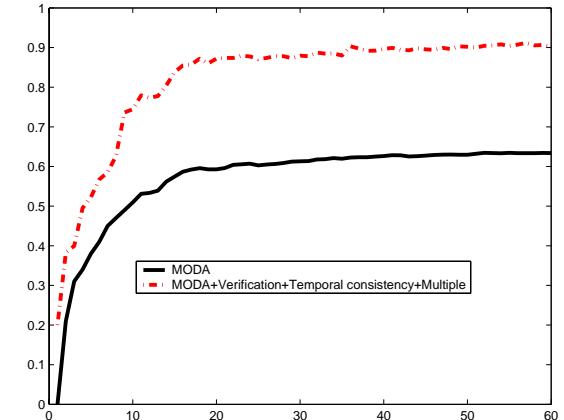
### Challenges

- Low quality small images (40-50 pixels).
- Changes in expression/pose/occlusion/illumination.
- Real time and scalable to several users.



## Adding space-time constraints

(de la Torre et al., 2005b)



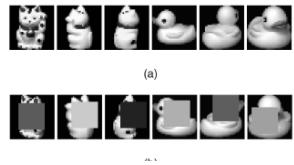
## 10- Robust Linear Discriminant Analysis

(Fidler et al., 2006)

- How to construct robust discriminative models?

$$\text{No noise model.}$$

$$G = B^T D$$



- Remove outliers in the space of generative models and learn linear a mapping to discriminative models.
- Assume ( $d > n$ ) and all the samples independent.

$$D = B^g C \quad \text{generative}$$

$$G = B^{d^T} D \quad \text{discriminative}$$

$$D \in \mathbb{R}^{d \times n} \quad B^g \in \mathbb{R}^{d \times n} \quad C \in \mathbb{R}^{n \times n}$$

$$B^d \in \mathbb{R}^{n \times n} \quad G \in \mathbb{R}^{n \times c}$$

samples  
pixels  
classes



Component Analysis for Computer Vision

F. De la Torre

ECCV-06

81

## Robust LDA

- $B^d$  as linear combination of the basis of  $B^g$ .

$$B^d = B^g V = \left[ B_{1:k}^g \ B_{n-k}^g \begin{bmatrix} V_{1:k} \\ V_{n-k} \end{bmatrix} \right] \quad V \in \mathbb{R}^{n \times c}$$

PCA coefficients

$$G = B^{d^T} D = V^T B^g T D = V^T C = V^T C_{1:k} + V^T C_{n-k:n}$$

- If just  $V_{1:k}^T C_{1:k}$  discriminative power can be lost.
- Add  $c$  basis to fully recover the LDA solution.

$$W = B_{n-k}^g V_{n-k:n} \in \mathbb{R}^{d \times c} \quad \hat{W} = W(W^T W)^{-\frac{1}{2}}$$

$$\hat{B}^s = [B_{1:k}^g \ \hat{W}] \quad \hat{V} = \begin{bmatrix} V_{1:k} \\ (W^T W)^{-\frac{1}{2}} \end{bmatrix}$$

No loss of information

Component Analysis for Computer Vision

F. De la Torre

ECCV-06

82

## Robust PCA

- Given training data find  $B^d$ ,  $V$ ,  $\hat{B}^s = [B_{1:k}^g \ \hat{W}]$
- Find robustly the coefficients in the  $\hat{B}^s$  basis such that:  $d_i \approx \hat{B}^s c_i$  and then estimate  $G$ .
- Robust estimation is achieved by hypothesize-and-test paradigm.

$$\begin{bmatrix} d_{1i} \\ d_{di} \end{bmatrix} \approx \begin{bmatrix} b_{11}c_1 + \dots + b_{ik}c_k \\ b_{d1}c_d + \dots + b_{dk}c_k \end{bmatrix}$$



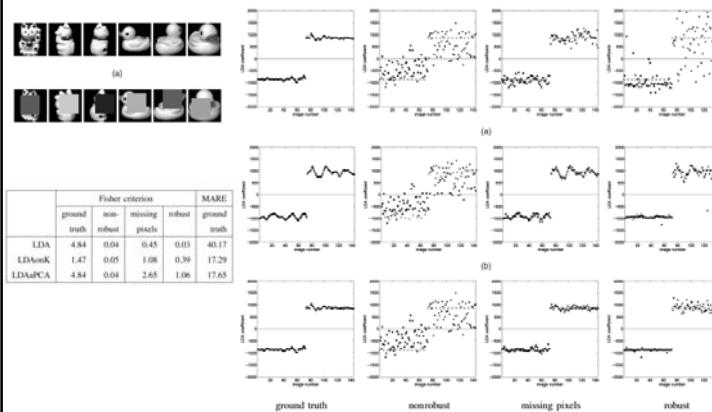
Component Analysis for Computer Vision

F. De la Torre

ECCV-06

83

## Experimental results



Component Analysis for Computer Vision

F. De la Torre

ECCV-06

84

## More examples

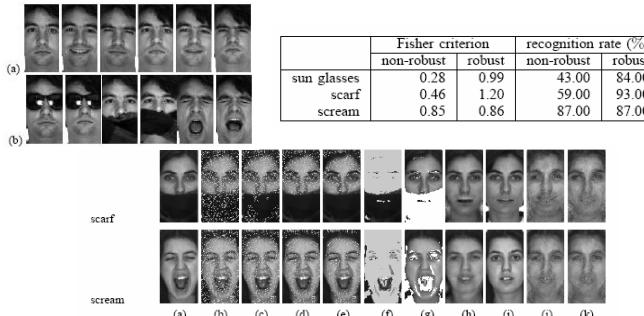


Fig. 8. Results of the pixel selection procedure. (a) test image, (b) random initializations in the pixel selection process (light-gray pixels denote the selected pixels), (c)-(e) a few stages of the  $\alpha$ -trimming method, (f) all compatible points, (g) the finally selected pixels used for coefficient calculation (denoted with the gray-scale of the original test image), (h) reconstructed image, (i) same image of the person in (h), (j) reconstructed image obtained with the standard LDA method, (k) reconstructed image obtained using naive ( $k = 0$ ) robust LDA method (see text for details).

## 11-Discriminative Cluster Analysis (DCA)

(de la Torre & Kanade, 2006)

- Generative clustering (e.g. k-means):

$$E(\mathbf{G}, \mathbf{B}) = \| \mathbf{D} - \mathbf{B}\mathbf{G}^T \|_F = \sum_{i=1}^c \sum_{j \in C_i} \| \mathbf{d}_j - \mathbf{b}_i \|$$

$$g_{ij} \in \{0,1\} \quad \mathbf{G}\mathbf{1}_k = \mathbf{1}_n$$

- Not efficient for high dimensional data.
- Multiple local minima.

- Discriminative clustering (de la Torre & Kanade, 2006):

$$E(\mathbf{V}, \mathbf{B}, \mathbf{G}) = \| (\mathbf{G}^T \mathbf{G})^{-\frac{1}{2}} (\mathbf{G}^T - \mathbf{V}\mathbf{B}^T \mathbf{D}) \|_F$$

- Simultaneous dimensionality reduction and clustering.

## Optimization

- Eliminate  $\mathbf{V}$

$$E(\mathbf{B}, \mathbf{G}) \propto \text{tr}((\mathbf{B}(\mathbf{D}\mathbf{D}^T\mathbf{B})^{-1}(\mathbf{B}^T\mathbf{D}\mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T\mathbf{D}^T\mathbf{B}))$$

- Optimize for  $\mathbf{B}$

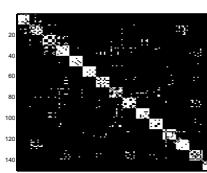
$$\mathbf{D}\mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T\mathbf{D}^T\mathbf{B} = \mathbf{D}\mathbf{D}^T\mathbf{B}\Lambda$$

- Optimize for  $\mathbf{G}$      $\mathbf{A} = \mathbf{C}^T(\mathbf{C}^T\mathbf{C})^{-1}\mathbf{C}$      $\mathbf{C} = \mathbf{B}^T\mathbf{D}$

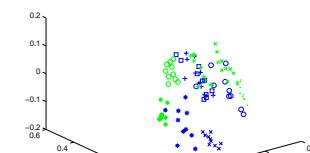
$$\mathbf{G} = \mathbf{V} \circ \mathbf{V} \quad \mathbf{V}^{(n+1)} = \mathbf{V}^{(n)} - \eta \frac{\partial E}{\partial \mathbf{V}}$$

$$\frac{\partial E}{\partial \mathbf{V}} = (\mathbf{I}_C - \mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T)\mathbf{A}\mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}$$

## Experiments

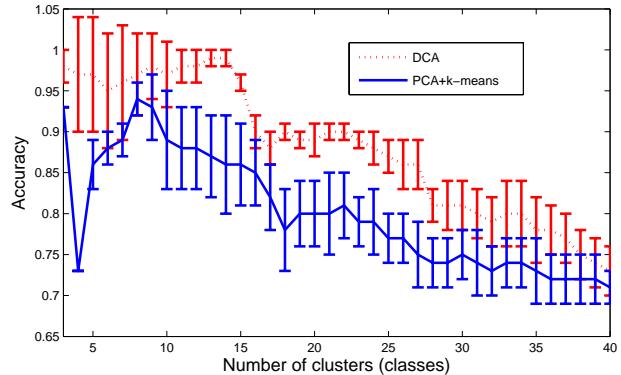


PCA

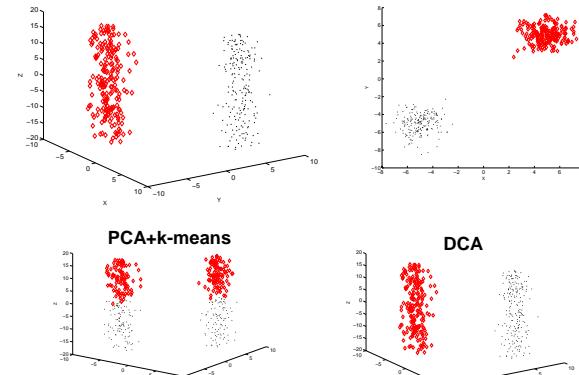


DCA

## DCA vs. PCA+k-means



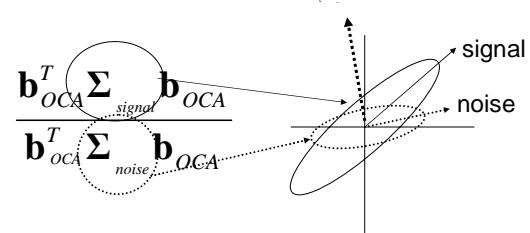
## Toy problem



## Related LDA work

- Face recognition (Belhumeur et al., 1997; Zhao, 2000; Martinez & Kak, 2003)
- Small sample problem (Chen et al., 2000; Yu & Yang, 2001)
- Mixture (Hastie et al., 1995; Zhu & Martinez, 2006;)
- Neural approaches (Gallinari et al., 1991; Lowe & Webb, 1991)
- Heteroscedastic Discriminant Analysis (Kumar & Andreou, 1998; Fukunaga, 1990; Mardia et al., 1979; Saon et al., 2000;)

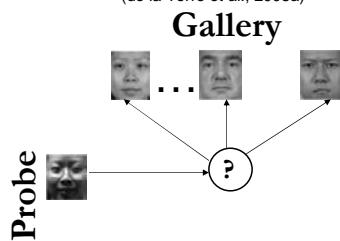
## Oriented Component Analysis (OCA)



- Generalized eigenvalue problem:  $\Sigma_i \mathbf{b}_k = \Sigma_e \mathbf{b}_k \lambda$
- $\mathbf{b}_{oca}$  is steered by the distribution of noise.

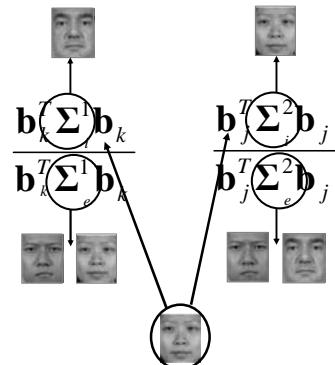
## 12- Representational Oriented Component Analysis (ROCA)

(de la Torre et al., 2005a)



- Challenges
  - Just 1 training image.
  - Changes in appearance, expression and illumination.

## OCA



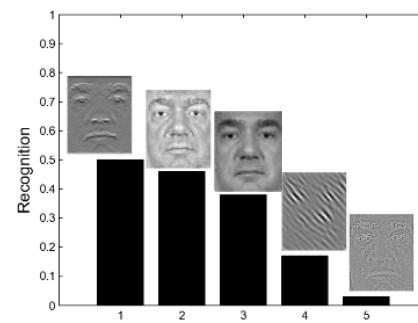
## Interpretation of OCA

(de la Torre et al., 2005a)

$$\begin{aligned} & \text{Probe Image} \\ & \xrightarrow{\mathbf{b}_k^T \Sigma_e^{-1} \mathbf{b}_k} \mathbf{b}_k = \Sigma_e^{-1} \mathbf{d}_k \\ & \quad \Sigma_e \approx \mathbf{U} \Lambda \mathbf{U}^T + \sigma^2 \mathbf{I} \\ & \mathbf{b}_k = \Sigma_e^{-1} \mathbf{d}_k \propto (\mathbf{I} - \mathbf{U}) \begin{pmatrix} \frac{\lambda_1 - \sigma^2}{\lambda_1} & & \\ & \frac{\lambda_2 - \sigma^2}{\lambda_2} & \\ & & \ddots \end{pmatrix} \mathbf{U}^T \mathbf{d}_k \end{aligned}$$

## Does representation matter?

- Unstable classifier with large probability of misclassification (prone to over-fitting).



## Combining several representations

Gallery →

- Morphological filters, Gabor filters, edge detectors, derivatives of Gaussians, etc.. (150 images)



$$\max_{\mathbf{B}_k^i} \frac{\left| \mathbf{B}_k^{i T} \left[ \mathbf{B}_k^i + \mathbf{B}_k^i (\mathbf{B}_k^i)^T + (\mathbf{B}_k^i - \mathbf{B}_k^i) (\mathbf{B}_k^i - \mathbf{B}_k^i)^T \right] \mathbf{B}_k^i \right|}{\left| \mathbf{B}_k^{i T} \left[ \dots + \mathbf{B}_k^i + \mathbf{B}_k^i + \dots \right] \mathbf{B}_k^i \right|}$$

## Solving Generalized Eigenvalue

- Rank deficient matrices.  $\|\Sigma_i - \mathbf{U}_i \mathbf{U}_i^T - \sigma^2 \mathbf{I}\|_F$   $\Sigma_i \approx \mathbf{U}_i \mathbf{U}_i^T + \sigma^2 \mathbf{I}$
- Not numerically stable algorithms/bad generalization (over fitting)

$$\max_{\mathbf{B}_k^i} \frac{\left| \mathbf{B}_k^{i T} (\mathbf{A} \mathbf{B}_k^i) \right|}{\left| \mathbf{B}_k^{i T} (\mathbf{C} \mathbf{B}_k^i) \right|} \rightarrow 10.000 \times 10.000$$

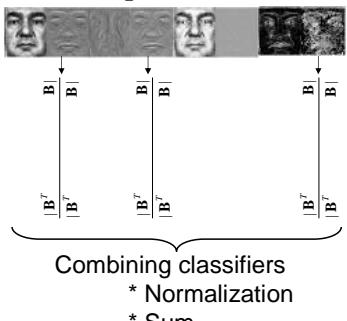
100      100  
100      100

- Modified Subspace Iteration.

$$\begin{aligned} \hat{\mathbf{C}}\hat{\mathbf{V}}_{k+1} &= \mathbf{A}\hat{\mathbf{V}}_k & \hat{\mathbf{V}}_{k+1} &= \hat{\mathbf{V}}_{k+1}/\max(\hat{\mathbf{V}}_{k+1}) \\ \mathbf{S} &= \hat{\mathbf{V}}_{k+1}^T \mathbf{A} \hat{\mathbf{V}}_{k+1} & \mathbf{T} &= \hat{\mathbf{V}}_{k+1}^T \mathbf{C} \hat{\mathbf{V}}_{k+1} \\ \mathbf{SW} &= \mathbf{CW}\Delta & & \\ \mathbf{V}_{k+1} &= \hat{\mathbf{V}}_{k+1} \mathbf{W} & & \end{aligned}$$

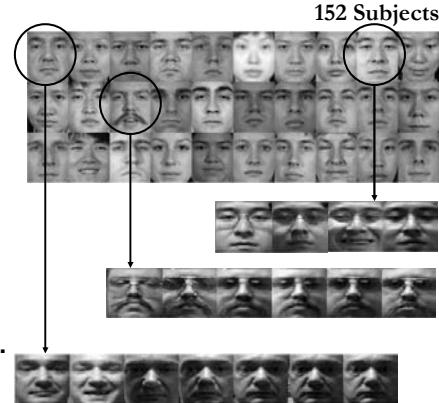
## Combining Classifiers

150 Representations

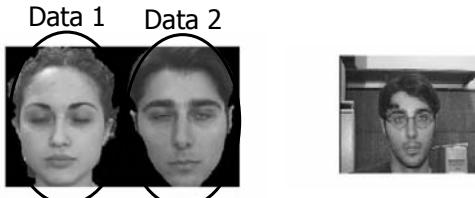


## Results

- PCA –
  - Mah 10%
  - Euclidean 16%
- NN – 29%
- Face it – 40.9%
- ROCA
  - Max of individual classifiers 47%.
  - Combining 62.1%.



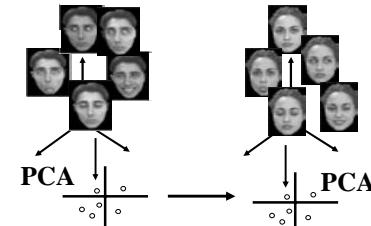
## 13- Dynamic Coupled Component Analysis (DCCA) (de la Torre & Black, 2001a)



- Learning the coupling.
- High dimensional data.
- Limited training data.

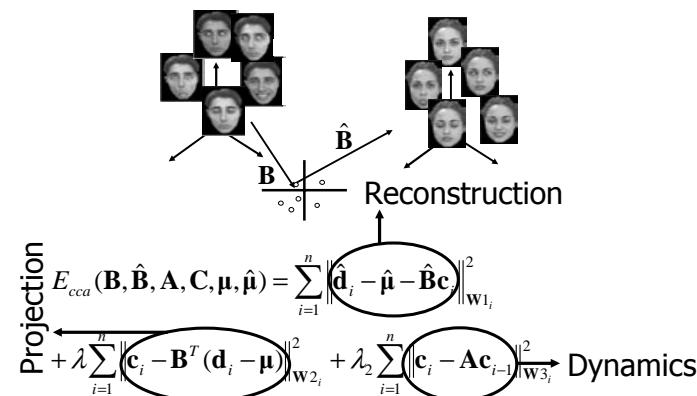
## Solutions?

- PCA independently and general mapping

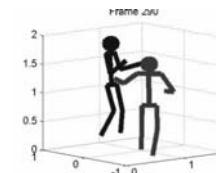


- Signals dependent signals with small energy can be lost.

## DCCA

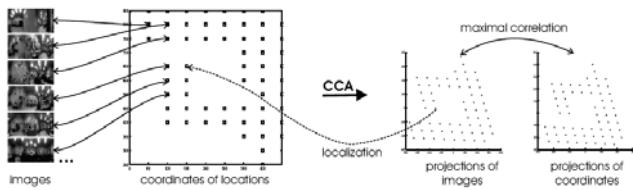


## Dynamic Coupled Component Analysis



## Robot localization with Canonical Correlation Analysis

(Skocaj & Leonardis, 2000)



Component Analysis for Computer Vision

F. De la Torre

ECCV-06

105

## Canonical Correlation Analysis (CCA)

(Mardia et al., 1979; Borga)

- Learn relations between multiple data sets? (e.g. find features in one set related to another data set)
- Given two sets  $\mathbf{X} \in \Re^{d_1 \times n}$  and  $\mathbf{Y} \in \Re^{d_2 \times n}$ , CCA finds the pair of directions  $\mathbf{w}_x$  and  $\mathbf{w}_y$  that maximize the correlation between the projections (assume zero mean data)

$$\rho = \frac{\mathbf{w}_x^T \mathbf{X}^T \mathbf{Y} \mathbf{w}_y}{\sqrt{\mathbf{w}_x^T \mathbf{X}^T \mathbf{X} \mathbf{w}_x \mathbf{w}_y^T \mathbf{Y}^T \mathbf{Y} \mathbf{w}_y}}$$

- Several ways of optimizing it:

$$\mathbf{A} = \begin{bmatrix} \mathbf{0} & \mathbf{X}^T \mathbf{Y} \\ \mathbf{X}^T \mathbf{Y} & \mathbf{0} \end{bmatrix} \in \Re^{(d_1+d_2) \times (d_1+d_2)}, \quad \mathbf{B} = \begin{bmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{0} \\ \mathbf{0} & \mathbf{Y}^T \mathbf{Y} \end{bmatrix} \in \Re^{(d_1+d_2) \times (d_1+d_2)} \quad \mathbf{w} = \begin{bmatrix} \mathbf{w}_x \\ \mathbf{w}_y \end{bmatrix}$$

- An stationary point of  $r$  is the solution to CCA.

$$\mathbf{A}\mathbf{w} = \lambda \mathbf{B}\mathbf{w}$$

Component Analysis for Computer Vision

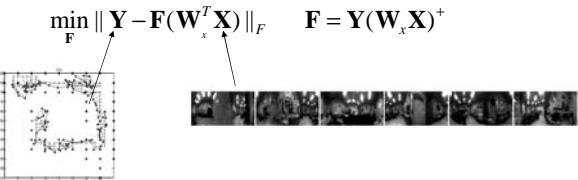
F. De la Torre

ECCV-06

106

## More on Canonical Correlation Analysis

- If  $d_1 >> n$  using the kernel trick efficient ways of solving it.
- Maximum number of canonical correlation vectors  $\min(d_1, d_2)$
- Learn a linear mapping between  $\mathbf{Y}$  and the projection of  $\mathbf{X}$  into canonical components.



Component Analysis for Computer Vision

F. De la Torre

ECCV-06

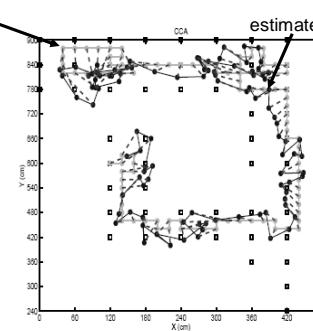
107

## Examples

### • Training images



True



(a) 2DOF	k	error	secs	elements
CCA	2	34.77	0.02	3906
PCA	2	77.37	0.21	4148
PCA	10	26.16	0.23	15044
PCAAir	10	13.04	2.77	24500
KCCA	10	26.65	0.23	15044
KCCAAir	10	13.12	2.81	24500
CCAAir	2	58.24	0.02	3906
CCAAcos	2	34.33	0.02	3906
CCAtan	4	30.65	0.02	6520
CCAsinccos	4	27.32	0.03	6520

Component Analysis for Computer Vision

F. De la Torre

ECCV-06

108

## Relevant Component Analysis

(Shental et al., 2002)

- Adding side-information



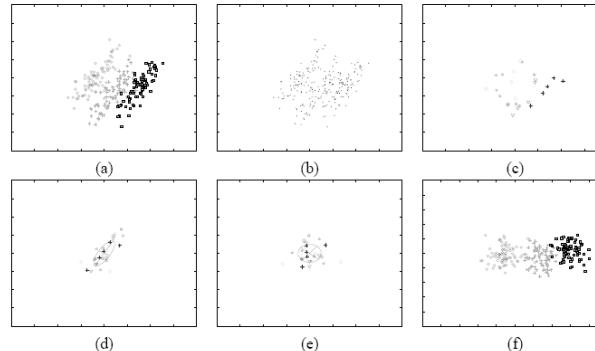
### Algorithm 1 The RCA algorithm

Given a dataset  $\{x_i\}_{i=1}^N$  and  $k$  chunklets  $C_j = \{x_{ji}\}_{i=1}^{n_j}$ ,  $j = 1 \dots k$ , do

- For each chunklet  $C_j$ , subtract the chunklet's mean from all the points it contains (Fig. 1d).
  - Compute the covariance matrix of all the centered data-points in chunklets (Fig. 1d). Assume a total of  $p$  points in  $k$  chunklets, where chunklet  $C_j$  consists of points  $\{x_{ji}\}_{i=1}^{n_j}$  and its mean is  $\hat{m}_j$ . RCA computes the following matrix:
- $$\hat{C} = \frac{1}{p} \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ji} - \hat{m}_j)(x_{ji} - \hat{m}_j)^T \quad (1)$$
- Compute the whitening transformation  $W = \hat{C}^{-\frac{1}{2}}$  associated with this covariance matrix (Fig. 1e), and apply it to the original data points:  $x_{new} = Wx$  (Fig. 1f). Alternatively, use the inverse of  $\hat{C}$  as a Mahalanobis distance.



## Example



## Standard extensions

- Latent Variable Models
- Tensor Factorization
  - 2D PCA/LDA.
  - Higher order extension.
- Kernel Methods



## Factor Analysis

- A Gaussian distribution on the coefficients and noise is added to PCA → Factor Analysis. (Mardia et al., 1979)

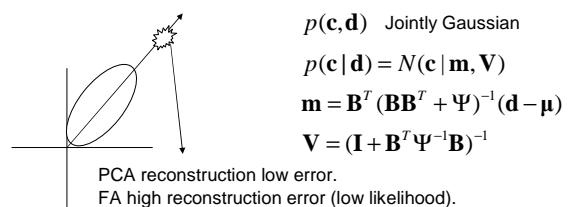
$$\mathbf{d} = \boldsymbol{\mu} + \mathbf{B}\mathbf{c} + \boldsymbol{\eta}$$

$$p(\mathbf{c}) = N(\mathbf{c} | \mathbf{0}, \mathbf{I}_k) \quad p(\mathbf{d} | \mathbf{c}, \mathbf{B}) = N(\mathbf{d} | \boldsymbol{\mu} + \mathbf{B}\mathbf{c}, \boldsymbol{\Psi})$$

$$p(\boldsymbol{\eta}) = N(\boldsymbol{\eta} | \mathbf{0}, \boldsymbol{\Psi}) \quad \boldsymbol{\Psi} = diag(\eta_1, \eta_2, \dots, \eta_d)$$

$$\text{cov}(\mathbf{d}) = E((\mathbf{d} - \boldsymbol{\mu})(\mathbf{d} - \boldsymbol{\mu})^T) = \mathbf{B}\mathbf{B}^T + \boldsymbol{\Psi}$$

- Inference (Roweis & Ghahramani, 1999; Tipping & Bishop, 1999a)

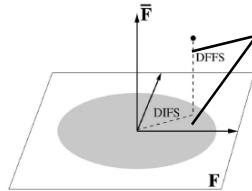


## Ppca

- If  $\Psi = E(\mathbf{p}\mathbf{p}^T) = \epsilon I_d$  PPCA.
- If  $\epsilon \rightarrow 0$  is equivalent to PCA.  $\epsilon \rightarrow 0$   $\mathbf{B}^T(\mathbf{B}\mathbf{B}^T + \Psi)^{-1} = (\mathbf{B}^T\mathbf{B})^{-1}\mathbf{B}^T$
- Probabilistic visual learning (Moghaddam & Pentland, 1997;)

$$p(\mathbf{d}) = \int p(\mathbf{d} | \mathbf{c}) p(\mathbf{c}) d\mathbf{c} = \frac{e^{-\frac{1}{2}(\mathbf{d}-\mu)^T \Sigma^{-1} (\mathbf{d}-\mu)}}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} = \frac{e^{-\frac{1}{2}(\mathbf{d}-\mu)^T (\mathbf{B}\mathbf{B}^T + \epsilon I_d)^{-1} (\mathbf{d}-\mu)}}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} = \frac{\left[ e^{-\frac{1}{2} \sum_{i=1}^k \frac{c_i^2}{\lambda_i}} \right] \left[ \frac{\epsilon^{\frac{d}{2}} (\mathbf{d})}{(2\pi)^{\frac{d}{2}} \prod_{i=1}^k \lambda_i^{1/2}} \right]}{(2\pi)^{\frac{d}{2}} (2\pi\mu)^{\frac{(d-k)}{2}}}$$

$$\mathbf{c}_i = \mathbf{B}^T \mathbf{d}_i$$



## A least squares interpretation

- Directly minimizing

$$E(\mathbf{B}, \Lambda, \sigma) = \|\Sigma - \mathbf{B}\Lambda\mathbf{B}^T - \sigma^2 \mathbf{I}_d\|_F$$

$\mathbf{B}^T\mathbf{B} = \mathbf{I}_d$   $\Lambda$  is diagonal

derives in the same solution as PPCA (de la Torre & Kanade, 2005b).

$$\begin{aligned} \sum \text{d} \times \text{d} & \quad \mathbf{B}\mathbf{B}^T + \sigma^2 \mathbf{I} \\ \frac{d(d+1)}{2} & \approx \frac{d^2}{2} \quad \frac{k(2d-k+1)}{2} \approx kd \end{aligned}$$

## More on PPCA

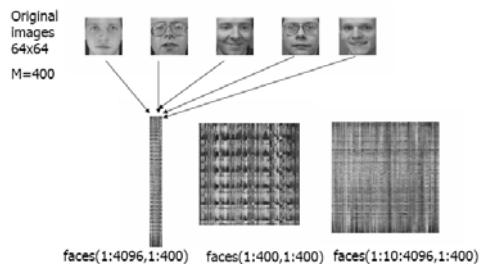
- Extension to mixtures of Ppca (mixture of subspaces). (Tipping & Bishop, 1999b; Black et al., 1998; Jebara et al., 1998)
- Tracking (Yang et al., 1999; Yang et al., 2000a; Lee et al., 2005; de la Torre et al., 2000b)
- Recognition/Detection (Moghaddam et al., 2000; Shakhnarovich & Moghaddam, 2004; Everingham & Zisserman, 2006)
- PCA for the exponential family (collins et al., 2001)

## Tensor Decomposition

- 2D PCA/LDA
- General tensor factorization.

## 2D PCA/SVD

- Vectorizing images do not preserve 2D properties and spatial properties are lost.
- Many ways of stacking images into matrices



## SVD versus 2D SVD



(Ye, 2004, Ding & Ye, 2006)

$$E(\mathbf{L}, \mathbf{R}, \{\mathbf{M}_i\}) = \sum_{i=1}^n \left\| \mathbf{D}_i^T - \mathbf{L}\mathbf{M}_i\mathbf{R}^T \right\|_F$$

$$\mathbf{D}_i \in \Re^{r \times c} \quad \mathbf{L} \in \Re^{r \times l_1} \quad \mathbf{R} \in \Re^{c \times l_2} \quad \mathbf{M}_i \in \Re^{l_1 \times l_2}$$

$$\mathbf{L}^T \mathbf{L} = \mathbf{I} \quad \mathbf{R}^T \mathbf{R} = \mathbf{I}$$

- Compression ratio.  $\frac{nrc}{rl_1 + cl_2 + nl_2l_1}$
- Recognition.  $\|\mathbf{D}_i - \mathbf{D}_j\| = \|\mathbf{L}(\mathbf{M}_i - \mathbf{M}_j)\mathbf{R}\| \approx \|\mathbf{(M}_i - \mathbf{M}_j)\|$
- 2D SVD smaller computational cost (space & time)
- Same or better reconstruction (for same number of parameters)

$$\frac{rc}{l_2l_1}$$

## Optimization

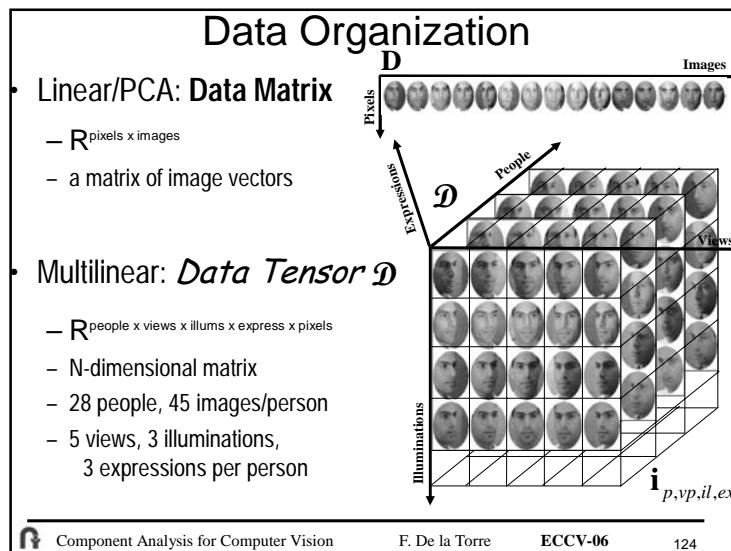
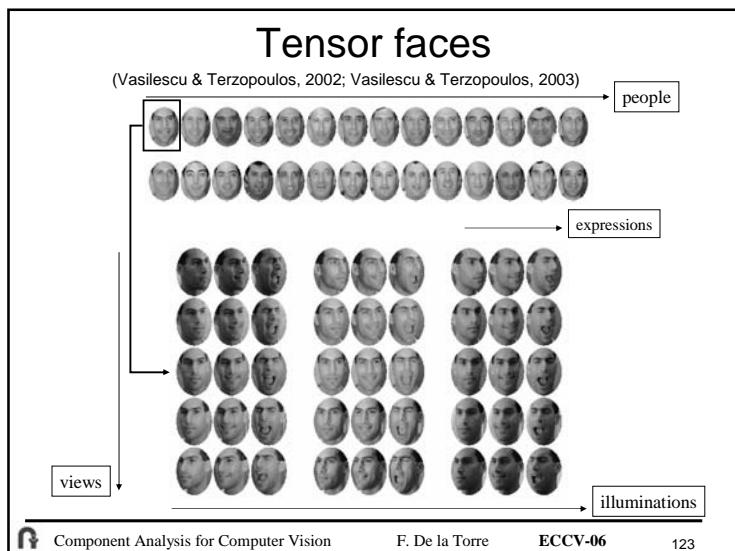
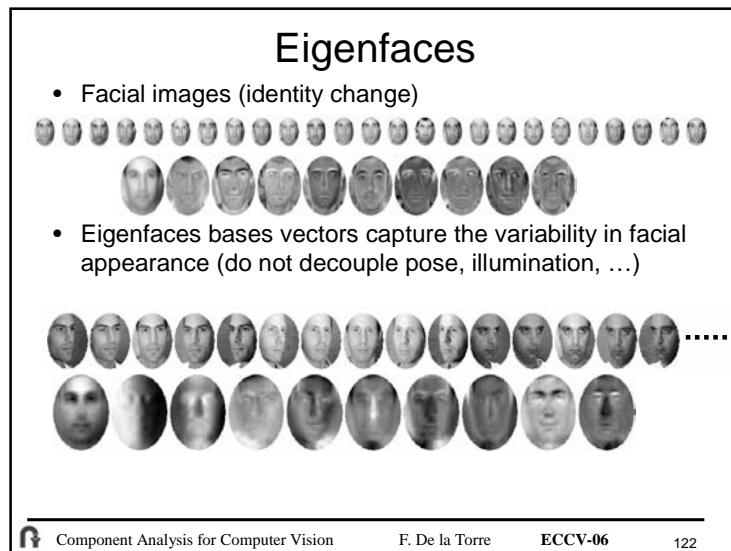
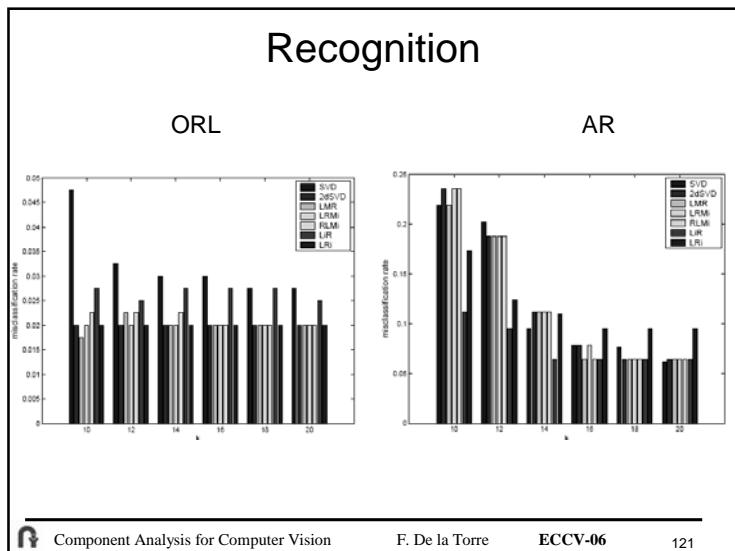
- No closed form solution.
- Alternate

Iterate until convergence 
$$\begin{cases} \left( \sum_{i=1}^n \mathbf{D}_i \mathbf{R} \mathbf{R}^T \mathbf{D}_i^T \right) \mathbf{L} = \mathbf{L} \Lambda_1 \\ \left( \sum_{i=1}^n \mathbf{D}_i^T \mathbf{L} \mathbf{L}^T \mathbf{D}_i \right) \mathbf{R} = \mathbf{R} \Lambda_2 \\ \mathbf{M}_i = \mathbf{L}^T \mathbf{D}_i \mathbf{R} \end{cases}$$
 Iterate until convergence

## 2D PCA/SVD

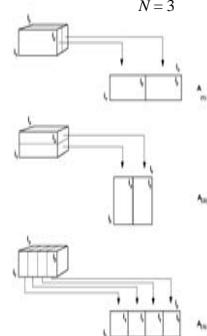
- SVD (k=15), storage 160560
- 2DSVD ( $l_1=l_2=15$ ), storage 93060





## N-Mode SVD Algorithm

$$D = Z \times_1 U_{\text{people}} \times_2 U_{\text{views}} \times_3 U_{\text{illums.}} \times_4 U_{\text{express.}} \times_5 U_{\text{pixels}}$$



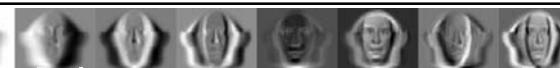
Component Analysis for Computer Vision

F. De la Torre

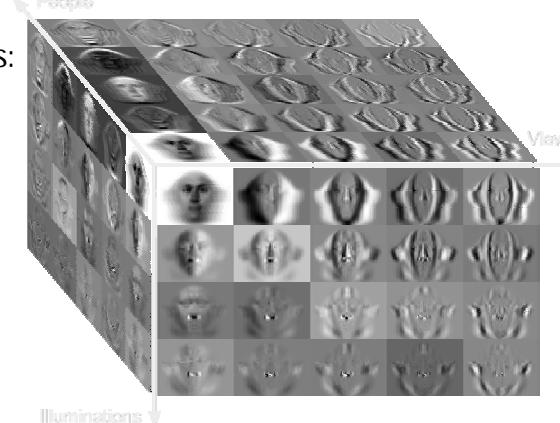
ECCV-06

125

PCA:



TensorFaces:



Component Analysis for Computer Vision

F. De la Torre

ECCV-06

126

## Strategic Data Compression = Perceptual Quality

- TensorFaces data reduction in illumination space primarily degrades illumination effects (cast shadows, highlights)
- PCA has *lower mean square error* but *higher perceptual error*

	TensorFaces	TensorFaces	PCA
Original	6 illum + 11 people param.	Mean Sq. Err. = <b>409.15</b>	Mean Sq. Err. = <b>85.75</b>
176 basis vectors	66 basis vectors	3 illum + 11 people param.	33 parameters

Component Analysis for Computer Vision

F. De la Torre

ECCV-06

127

## Results

Data Set - 16,875 images

- 75 people
- 15 viewpoints
- 15 illuminations

Training Images - 2,700

- 75 people
- 6 viewpoints
- 6 illuminations

Test Images:

- 75 people
- 9 viewpoints
- 9 illums

Linear Models		Multilinear Models	
PCA	ICA	TensorFaces	Independent TensorFaces
83%	89%	93%	97%

Component Analysis for Computer Vision

F. De la Torre

ECCV-06

128

## Related work

- **Tensor factorization** (O'Leary & Peleg, 1983; Shashua & Levin, 2001; Paatero & Tapper, 1994; Shashua & Hazan, 2005)
- **2D PCA** (Kong & Wang, 2005; Ding & Ye, 2006; Zhang et al., 2006; Zhang & Zhou, 2005; Yang et al., 2004b)
- **2D LDA** (Ye, 2005; Ye et al., 2005; Liu et al., 1993)



## Kernel Methods

- Given in the tutorial.

